

# The Consumer Digital Infrastructure and Scientific Computing

Dr. David P. Anderson  
Space Sciences Laboratory, U.C. Berkeley  
Sept. 15, 2010

**Summary:** the consumer digital infrastructure has the potential to increase the amount of computing power available to U.S. scientists by orders of magnitude, with no hardware investment by funding agencies. The technology to realize this potential is already in place; what is needed is planning, organization and a marketing strategy.

## 1. The consumer digital infrastructure

Today's consumer digital infrastructure (CDI) consists of mass-market devices (e.g., desktop, laptop and tablet computers, game consoles, PDAs, media players, smart phones, VTRs) and the communication networks that connect them (the commodity Internet, cell phone and other radio networks). Driven by the requirements of Internet streaming video and 3-D graphical games, consumer products have become very powerful. Today's PC has 4 GB RAM, 1 TB disk, and a 1 TeraFLOPS GPU; it is capable of running almost any scientific application. Home network connections are on the order of 10 Mbps, increasing soon to 100 Mbps and then to the Gbps range as optical fiber to homes is deployed.

The CDI currently includes over 1 billion privately-owned PCs and 100 million GPUs capable of general-purpose computing. These have a total computing capability of roughly 100 ExaFLOPS. They have on the order of 10 Exabytes of free disk space, accessible via 1 Peta-bps of network bandwidth. In order to handle peak loads, the CDI is over-provisioned and hence underutilized. The average CPU and home Internet connection are used only a few percent of the time they are available – which, even with power-saving modes, is a significant fraction of the time.

The CDI is a possible platform for high-performance scientific computing, alongside platforms such as supercomputers, clusters, grids, and clouds. Each of these platforms has its strengths and limitations; in particular, the CDI has several advantages relative to the other platforms:

- The CDI has much larger processing and storage capacity, and hence enables otherwise infeasible science. For example, the Square Kilometer Array radio telescope will generate 0.1 to 1 TB/sec of data. Using the CDI, data could be stored longer (months versus hours) and more science could be extracted from it.
- The CDI's HVAC, electrical power and hardware are paid for by consumers, and the hardware is continuously upgraded to state-of-the-art components. The value of the CDI's PCs alone is roughly \$1 trillion.
- The CDI is self-maintaining: consumers fix their own software and hardware problems.
- Wide geographic, political and infrastructural distribution increases system robustness by eliminating single points of failure.
- Consumer products are the main focus of computing research and development, and have a substantial price/performance advantage over data-center hardware.

## 2. Volunteer computing and BOINC

To use the excess capacity of the CDI for scientific computing, we first need the consent of the resource owners. This requires incentives that reward consumers for the use of their computing resources, and publicity that makes the public aware of this possibility. Experience has shown that consumers can be motivated by their support for the goals of the research, by the chance to participate in online communities, and by competition

based on computational contribution. This is called **volunteer computing**.

Second, we need a software system to dispatch jobs from scientists' servers to consumer computers, and to execute jobs on those computers. This system must handle a variety of factors that are unique to volunteer computing:

- Scale: millions of nodes and tens of millions of jobs per day.
- Heterogeneity: substantial diversity of software, hardware, and network connectivity.
- Trust and reliability: incorrect computational results may be intentionally returned by malicious participants.
- Communication access: many systems are behind firewalls that allow only outgoing HTTP traffic.
- Availability: sporadic presence and a high churn rate.
- Ease of use: the client software must be extremely simple to install and must work with no configuration or user intervention.

BOINC [1], an NSF-funded project at UC Berkeley, has developed software addressing these issues. BOINC provides server software that lets scientists create volunteer computing projects, and client software, available for all major platforms, that lets volunteers participate in any combination of these projects. BOINC addresses the needs of volunteer computing; Grid platforms do not.

BOINC allows scientists to create **projects**, each of which consists of a web site and a server that dispatches jobs. Volunteers can **attach** their computers to one or more projects, and can control the resource allocation among them. Projects may offer multiple applications, and a project may optionally allow volunteers to select which applications they wish to run. The BOINC architecture supports **account managers**: web services that provide a level of indirection between volunteer and projects. A computer may be attached to an account manager rather than to individual projects. The account manager periodically tells the computer which projects to do work for, and the computer gets jobs directly from these projects.

BOINC was released in 2004, and there are now about 50 volunteer computing projects and 3 account managers. The volunteer population consists of 500,000 people and 1 million computers, providing an average throughput of 12 PetaFLOPS. This has resulted in a large amount of research and publications, including papers in top journals like Nature, Science, Physics Review, and Cell [2].

### 3. Scientific applications of volunteer computing

Volunteer computing is useful for most “bag of tasks” applications: parameter sweeps, simulations with perturbed initial conditions, compute-intensive data analysis, genetic algorithms. BOINC is used by projects from many institutions, doing research in many areas, including astrophysics, cosmology, climate study, biochemistry, epidemiology, environmental science, cognitive science, genetics, mathematics, nanotechnology, particle physics, quantum computing, and seismology. BOINC supports many types of scientific applications:

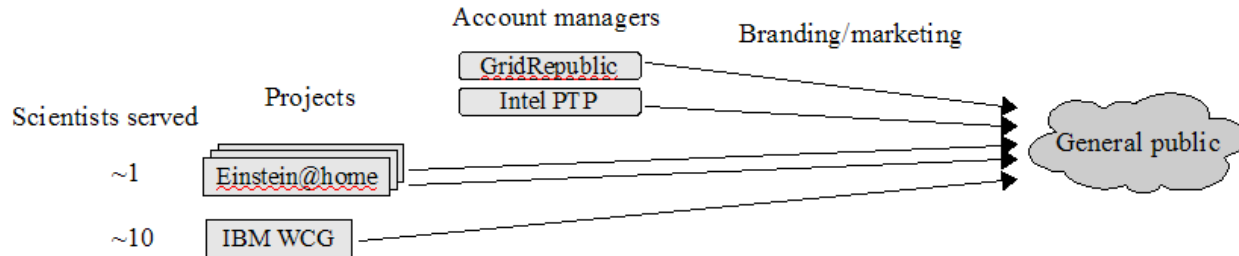
- Data-intensive applications; for example, the Einstein@home project analyzes Petabyte-scale data from the LIGO gravitational-wave observatory and the Arecibo radio observatory [3]. BOINC includes a data management system that optimizes network traffic for data-intensive applications.
- Applications in FORTRAN, C, Java, Python, CUDA, OpenCL, or other language systems.
- Multi-thread applications.
- Applications that run in virtual machines.
- Legacy applications for which only an executable is available.
- Jobs that take seconds, days, or months of processing time.
- Jobs with extreme memory, disk, and/or latency requirements; such jobs are sent only to hosts that are able to handle them.

BOINC does not currently support MPI-type applications. However, GPU-based computing now offers an alternative paradigm that is preferable in many cases.

#### 4. Problems, goals and ways forward

Realizing the potential of volunteer computing involves two related goals: a) making volunteer computing available to all scientists, and b) growing the volunteer population from a half million people to tens or hundreds of millions.

Volunteer computing's status quo can be depicted as follows:



There are 50 or so projects. Most are operated by individual research groups; one (IBM World Community Grid [4]) is an “umbrella” project serving 10 or so scientists. There are three account managers: **BAM!**, **Gridrepublic**, and a Facebook-based Intel promotion, **Progress Thru Processors** [5, 6, 7], which provide a “one-stop shopping” interface for browsing and attaching to projects. The projects and account managers all have their own brands, which are marketed separately to the public.

##### 4.1 Increasing accessibility to scientists

Creating and operating a BOINC project requires significant skills and resources, which only a few research groups have. IBM WCG offers an alternative, but it has various limitations – for example, it supports only applications deemed by IBM to be “humanitarian”.

One way to break this bottleneck is to leverage existing organizations that provide HPC to scientists:

- Supercomputer centers.
- Science portals such as NanoHub.
- National grids such as TeraGrid and Open Science Grid.
- Campus-level computing centers.

I propose that these **HPC providers** create and operate BOINC-based volunteer computing projects. The organizations are well-situated to do so: they have the staff, skills and resources to support scientists and to operate servers, and an existing clientele of scientists. A BOINC-based back end can be added transparently to these scientists.

There are two incentives for HPC providers to do this. First, volunteer computing complements rather than replaces existing resources. Parallel programs need supercomputers and MPI programs need clusters. By offloading other jobs onto the CDI, the availability and effectiveness of these specialized resources will be increased. Second, offering a volunteer computing project will increase the public visibility of the organization. In some cases, the organization has a unique opportunity to easily market itself to a particular population; for example, a UCB campus-level project could market itself to UCB alumni, of which there are about 450,000, via its alumni newsletter).

## 4.2 Increasing volunteership

Attracting volunteers is a marketing exercise, involving various channels (mass media, web ads, social networks, and so on), the creation of “brands”, and the design of incentives. In the status quo, each project markets itself to the general public, resulting in a glut of brands and a diluted and confusing marketing message to the public.

I believe that effectively marketing volunteer computing requires a unified brand, ideally at the national level. I propose the creation of an organization – say, **ScienceUSA.org** – representing the totality of volunteer computing for research funded by the US government. PC owners wishing to volunteer would simply go to ScienceUSA.org, install the software, and register for general areas, like “disease research”, “environmental research”, and “space and physics”, or for specific sub-areas. The web site would provide science information and updates on all the research involved. The various areas would be mapped dynamically, according to a resource allocation policy, to specific projects and applications within projects. In concrete terms, ScienceUSA.org would be a BOINC account manager, not a BOINC project.

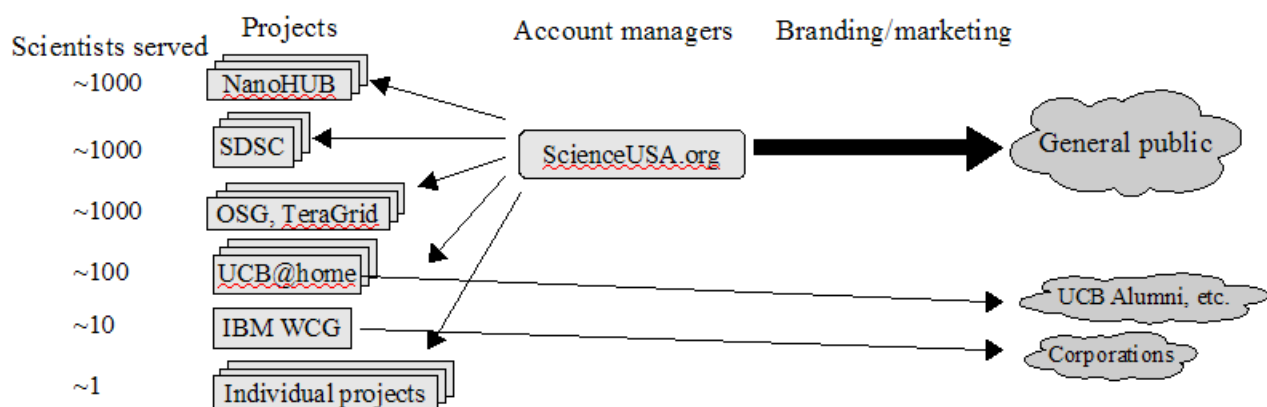
As an organization, ScienceUSA.org would perform the following tasks:

- Publicize itself (in national media, social networks, educational and professional channels).
- Work with major OS distributors (Microsoft, Apple, Ubuntu) and PC vendors to include volunteer computing as an opt-in part of their default installation. This is already being done on a small scale: volunteer computing software is bundled with Sony Playstation 3s and with ATI GPU drivers.
- Develop and maintain its web site.
- Coordinate with the various BOINC-based projects.
- Create a framework for defining and enforcing a resource allocation policy.

The resource allocation policy would be decided by a committee representing the funding agencies, HPC providers, and other stakeholders. Funding agency awards could include a fraction of the aggregate computing power; in some cases this could replace awards for HPC hardware purchases.

## 4.3 Implementing the ScienceUSA.org model

The model proposed above can be depicted as follows:



Developing this model involves the following steps:

- Form a **volunteer computing task force**, with representation from multiple funding agencies, from supercomputing centers, hubs and grids, and from scientific application areas. This task force would articulate the organizational structure proposed above, and propose a resource allocation policy.
- Create volunteer computing projects at HPC providers (supercomputing centers, hubs, and grids), and

- the develop software to interface their existing job-submission systems to these projects.
- Create and operate of ScienceUSA.org.

These are low-cost activities (a few FTEs and some low-end hardware), and would result in a computing resource that would be usable by almost all computational scientists, and whose capacity would otherwise cost billions of dollars per year. ScienceUSA.org would also provide a powerful education and public outreach channel.

## References

- [1] <http://boinc.berkeley.edu>
- [2] [http://boinc.berkeley.edu/wiki/Publications\\_by\\_BOINC\\_projects](http://boinc.berkeley.edu/wiki/Publications_by_BOINC_projects)
- [3] <http://einstein.phys.uwm.edu/>
- [4] <http://www.worldcommunitygrid.org/>
- [5] <http://bam.boincstats.com/>
- [6] <http://www.gridrepublic.org/>
- [7] <http://www.facebook.com/progressthruprocessors>